

---

## Fisher Auto-Encoders: Supplementary Materials

---

### 1 Proof of Theorem 1

$$\begin{aligned}
\phi^*, \eta^*, \theta^* &= \arg \min_{\phi, \eta, \theta} \mathbb{D}_{\nabla} [q_{\star, \phi}(\mathbf{x}, \mathbf{z}) || p_{\theta}(\mathbf{x}, \mathbf{z})] \\
&= \arg \min_{\phi, \eta, \theta} \mathbb{E}_{q_{\star, \phi}(\mathbf{x}, \mathbf{z})} \frac{1}{2} \|\nabla_{\mathbf{x}, \mathbf{z}} \log q_{\star, \phi}(\mathbf{x}, \mathbf{z}) - \nabla_{\mathbf{x}, \mathbf{z}} \log p_{\eta, \theta}(\mathbf{x}, \mathbf{z})\|^2 \\
&= \arg \min_{\phi, \eta, \theta} \mathbb{E}_{p_{\star}(\mathbf{x})} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \frac{1}{2} \|\nabla_{\mathbf{x}, \mathbf{z}} \log q_{\star, \phi}(\mathbf{x}, \mathbf{z}) - \nabla_{\mathbf{x}, \mathbf{z}} \log p_{\eta, \theta}(\mathbf{x}, \mathbf{z})\|^2 \\
&= \arg \min_{\phi, \eta, \theta} \mathbb{E}_{p_{\star}(\mathbf{x})} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \frac{1}{2} \|\nabla_{\mathbf{x}} \log p_{\star}(\mathbf{x}) + \nabla_{\mathbf{x}} \log q_{\phi}(\mathbf{z}|\mathbf{x}) - \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}|\mathbf{z})\|^2 \\
&\quad + \underbrace{\mathbb{E}_{p_{\star}(\mathbf{x})} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \frac{1}{2} \|\nabla_{\mathbf{z}} \log q_{\phi}(\mathbf{z}|\mathbf{x}) - \nabla_{\mathbf{z}} \log p_{\eta, \theta}(\mathbf{z}|\mathbf{x})\|^2}_{\mathbb{D}_{\nabla} [q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\eta, \theta}(\mathbf{z}|\mathbf{x})]} \\
&= \arg \min_{\phi, \eta, \theta} \mathbb{E}_{p_{\star}(\mathbf{x})} \mathbb{D}_{\nabla} [q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\eta, \theta}(\mathbf{z}|\mathbf{x})] + \mathbb{E}_{p_{\star}(\mathbf{x})} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \frac{1}{2} \|\nabla_{\mathbf{x}} \log p_{\star}(\mathbf{x})\|^2 \\
&\quad + \mathbb{E}_{p_{\star}(\mathbf{x})} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \nabla_{\mathbf{x}} \log p_{\star}(\mathbf{x})^{\top} \nabla_{\mathbf{x}} \log q_{\phi}(\mathbf{z}|\mathbf{x}) \\
&\quad + \mathbb{E}_{p_{\star}(\mathbf{x})} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \frac{1}{2} \|\nabla_{\mathbf{x}} \log q_{\phi}(\mathbf{z}|\mathbf{x})\|^2 - \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}|\mathbf{z})^{\top} \nabla_{\mathbf{x}} \log q_{\phi}(\mathbf{z}|\mathbf{x}) \\
&\quad + \mathbb{E}_{p_{\star}(\mathbf{x})} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \frac{1}{2} \|\nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}|\mathbf{z})\|^2 - \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}|\mathbf{z})^{\top} \nabla_{\mathbf{x}} \log p_{\star}(\mathbf{x})
\end{aligned}$$

Let's examine the inner-product terms:

$$\begin{aligned}
&\mathbb{E}_{p_{\star}(\mathbf{x})} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \nabla_{\mathbf{x}} \log p_{\star}(\mathbf{x})^{\top} \nabla_{\mathbf{x}} \log q_{\phi}(\mathbf{z}|\mathbf{x}) \\
&= \iint p_{\star}(\mathbf{x}) q_{\phi}(\mathbf{z}|\mathbf{x}) \nabla_{\mathbf{x}} \log p_{\star}(\mathbf{x})^{\top} \nabla_{\mathbf{x}} \log q_{\phi}(\mathbf{z}|\mathbf{x}) d\mathbf{z} d\mathbf{x} \\
&= \iint p_{\star}(\mathbf{x}) q_{\phi}(\mathbf{z}|\mathbf{x}) \nabla_{\mathbf{x}} \log p_{\star}(\mathbf{x})^{\top} \nabla_{\mathbf{x}} \log q_{\phi}(\mathbf{z}|\mathbf{x}) d\mathbf{x} d\mathbf{z} \\
&= \iint p_{\star}(\mathbf{x}) \nabla_{\mathbf{x}} \log p_{\star}(\mathbf{x})^{\top} \nabla_{\mathbf{x}} q_{\phi}(\mathbf{z}|\mathbf{x}) d\mathbf{x} d\mathbf{z} \\
&\stackrel{(a)}{=} - \iint p_{\star}(\mathbf{x}) q_{\phi}(\mathbf{z}|\mathbf{x}) \left[ \|\nabla_{\mathbf{x}} \log p_{\star}(\mathbf{x})\|^2 + \Delta_{\mathbf{x}} \log p_{\star}(\mathbf{x}) \right] d\mathbf{x} d\mathbf{z} \\
&= -\mathbb{E}_{p_{\star}(\mathbf{x})} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \|\nabla_{\mathbf{x}} \log p_{\star}(\mathbf{x})\|^2 + \Delta_{\mathbf{x}} \log p_{\star}(\mathbf{x}) \right],
\end{aligned}$$

where (a) is obtained by an integration by parts.

$$\begin{aligned}
 & \mathbb{E}_{p_*(\mathbf{x})} \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} - \nabla_{\mathbf{x}} \log p_\theta(\mathbf{x}|\mathbf{z})^\top \nabla_{\mathbf{x}} \log q_\phi(\mathbf{z}|\mathbf{x}) \\
 &= - \iint p_*(\mathbf{x}) q_\phi(\mathbf{z}|\mathbf{x}) \nabla_{\mathbf{x}} \log p_\theta(\mathbf{x}|\mathbf{z})^\top \nabla_{\mathbf{x}} \log q_\phi(\mathbf{z}|\mathbf{x}) d\mathbf{z} d\mathbf{x} \\
 &= - \iint p_*(\mathbf{x}) \nabla_{\mathbf{x}} \log p_\theta(\mathbf{x}|\mathbf{z})^\top \nabla_{\mathbf{x}} q_\phi(\mathbf{z}|\mathbf{x}) d\mathbf{x} d\mathbf{z} \\
 &\stackrel{(b)}{=} \iint p_*(\mathbf{x}) q_\phi(\mathbf{z}|\mathbf{x}) [\Delta_{\mathbf{x}} \log p_\theta(\mathbf{x}|\mathbf{z}) + \nabla_{\mathbf{x}} \log p_\theta(\mathbf{x}|\mathbf{z})^\top \nabla_{\mathbf{x}} \log p_*(\mathbf{x})] d\mathbf{x} d\mathbf{z} \\
 &= \mathbb{E}_{p_*(\mathbf{x})} \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\Delta_{\mathbf{x}} \log p_\theta(\mathbf{x}|\mathbf{z}) + \nabla_{\mathbf{x}} \log p_\theta(\mathbf{x}|\mathbf{z})^\top \nabla_{\mathbf{x}} \log p_*(\mathbf{x})],
 \end{aligned}$$

where (b) is again obtained by an integration by parts. Grouping all the terms together, we get

$$\begin{aligned}
 & \phi^*, \eta^*, \theta^* \\
 &= \arg \min_{\phi, \eta, \theta} \mathbb{E}_{p_*(\mathbf{x})} - s_\nabla [p_*(\mathbf{x})] + \mathbb{D}_\nabla [q_\phi(\mathbf{z}|\mathbf{x}) || p_{\eta, \theta}(\mathbf{z}|\mathbf{x})] + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} s_\nabla [p_\theta(\mathbf{x}|\mathbf{z})] \\
 & \quad + \frac{1}{2} \|\nabla_{\mathbf{x}} \log q_\phi(\mathbf{z}|\mathbf{x})\|^2.
 \end{aligned}$$

By noticing that  $\mathbb{E}_{p_*(\mathbf{x})} - s_\nabla [p_*(\mathbf{x})]$  is independent of the parameters  $\phi$ ,  $\eta$  and  $\theta$ , we conclude the proof of Theorem 1.

## 2 Robustness to binary masking noise

We extend the experiments to examine the robustness of the proposed Fisher AEs and consider another type of noise called binary masking noise which consists on setting the value of a randomly selected fraction  $\nu$  of input components to zero.

### 2.1 MNIST

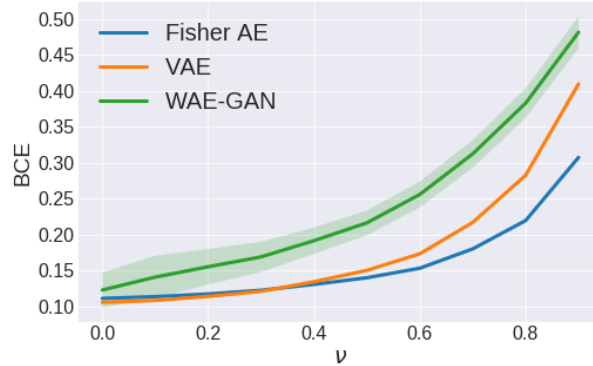


Figure 1: BCE vs. the fraction  $\nu$  for MNIST.

The same insights regarding the robustness of the Fisher AE to Gaussian noise are confirmed in the case of binary masking noise. Both Figures 1 and 2 shows the superiority of Fisher AE in terms of robustness to binary masking noise as compared to VAE and WAE-GAN.

### 2.2 CelebA

As shown in Figures 3 and 4, in the case of CelebA, both VAE and Fisher AE exhibit similar but superior performance in terms of robustness against binary masking noise as compared to WAE-GAN.

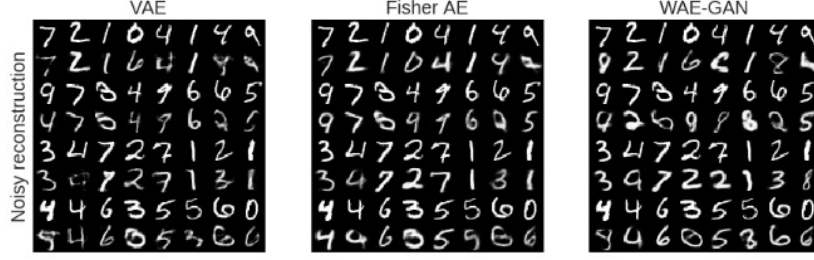


Figure 2: Test reconstruction results when a random fraction  $\nu = 0.8$  of test data is set to zero. True test data are given by the odd rows.

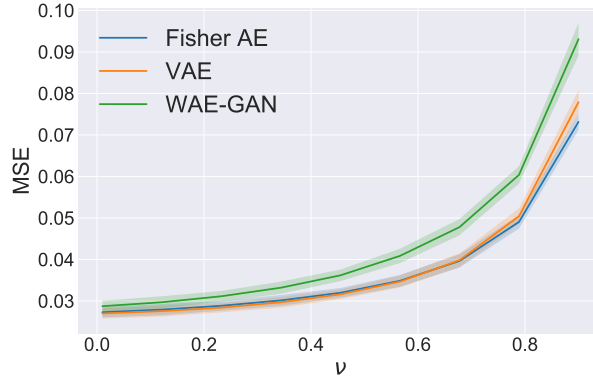


Figure 3: MSE vs. the fraction  $\nu$  for celebA with random mask.

### 3 Further details on experiments

Here, we give the detailed architecture used in the implementation of the different auto-encoders for both MNIST and celebA data sets.

- $\text{FC}(n_{in}, n_{out})$ : Fully connected layer with input/output dimensions given by  $n_{in}$  and  $n_{out}$ .
- $\text{Conv}(n_{in}, n_{out}, k, s, p)$ : Convolutional layer with input channels  $n_{in}$ , output channels  $n_{out}$ , kernel size  $k$ , stride  $s$  and padding  $p$ .
- $\text{ConvT}(n_{in}, n_{out}, k, s, p)$ : Transposed convolutional layer with input channels  $n_{in}$ , output channels  $n_{out}$ , kernel size  $k$ , stride  $s$  and padding  $p$ .
- $\text{AvgPool}(k, s, p)$ : Average Pooling with kernel size, stride and padding respectively given by  $k$ ,  $s$  and  $p$ .
- BN : Batch-normalization
- BiL : 2D bilinear interpolation layer



Figure 4: Test reconstruction results when a random fraction  $\nu = 0.8$  of test data is set to zero. True test data are given by the odd rows.

### 3.1 MNIST

Encoder	Decoder
Input size: (1, 28, 28) Conv(1, 64, 3, 2, 2) LeakyReLU Conv(64, 128, 3, 2, 2) BN LeakyReLU Conv(128, 256, 3, 2, 2) BN, LeakyReLU Conv(256, 512, 3, 2, 1) BN, LeakyReLU Conv(512, 16, 3, 2, 0) Output size : (16, 1, 1)	Input size: (8, 1, 1) ConvT(8, 512, 5, 1, 1) BN, ReLU ConvT(512, 256, 5, 1, 1) BN, ReLU ConvT(256, 128, 5, 2, 1) BN, ReLU ConvT(128, 64, 5, 1, 1) BN, ReLU ConvT(64, 1, 4, 2, 0) Sigmoid Output size : (1, 28, 28)

Table 1: Encoder/Decoder architectures for Fisher AE and VAE for MNIST

Encoder	Generator	Discriminator
Input size : (1, 28, 28) Conv(1, 64, 3, 2, 2) LeakyReLU Conv(64, 128, 3, 2, 2) BN, LeakyReLU Conv(128, 256, 3, 2, 2) BN, LeakyReLU Conv(256, 512, 3, 2, 1) BN, LeakyReLU Conv(512, 8, 3, 1, 0) Output size : (8, 1, 1)	Input size : (8, 1, 1) ConvT(8, 512, 5, 1, 1) BN, ReLU ConvT(512, 256, 5, 1, 1) BN, ReLU ConvT(256, 128, 5, 2, 1) BN, ReLU ConvT(128, 64, 5, 1, 1) BN, ReLU ConvT(64, 1, 4, 2, 0) Sigmoid Output size : (1, 28, 28)	Input size : (8, 1, 1) Flatten FC(8, 256) ReLU FC(256, 1) Sigmoid  Output size : (1, )

Table 2: Encoder/Generator/Discriminator architectures for WAE-GAN for MNIST

### 3.2 CelebA

Encoder	Decoder
Input size: (3, 64, 64) Conv(3, 64, 5, 1, 2) LeakyReLU, AvgPool(2, 2, 0) Conv(64, 128, 5, 1, 2) BN, LeakyReLU, AvgPool(2, 2, 0) Conv(128, 256, 5, 1, 2) BN, LeakyReLU, AvgPool(2, 2, 0) Conv(256, 512, 5, 1, 2) BN, LeakyReLU Flatten FC(8192, 64) Output size : (64, )	Input size: (64, 1, 1) BiL, Conv(64, 512, 5, 1, 02) BN, ReLU BiL, Conv(512, 256, 5, 1, 2) BN, ReLU BiL, Conv(256, 128, 5, 1, 2) BN, ReLU BiL, Conv(128, 64, 5, 1, 2) BN, ReLU BiL, Conv(64, 3, 5, 1, 2) Tanh Output size : (3, 64, 64)

Table 3: Encoder/Decoder architectures for Fisher AE and VAE for CelebA

Encoder	Generator	Discriminator
Input size: (3, 64, 64) Conv(3, 64, 5, 1, 2) LeakyReLU, AvgPool(2, 2, 0) Conv(64, 128, 5, 1, 2) BN, LeakyReLU, AvgPool(2, 2, 0) Conv(128, 256, 5, 1, 2) BN, LeakyReLU, AvgPool(2, 2, 0) Conv(256, 512, 5, 1, 2) BN, LeakyReLU Flatten FC(8192, 64) Output size : (64, )	Input size: (64, 1, 1) BiL, Conv(64, 512, 5, 1, 02) BN, ReLU BiL, Conv(512, 256, 5, 1, 2) BN, ReLU BiL, Conv(256, 128, 5, 1, 2) BN, ReLU BiL, Conv(128, 64, 5, 1, 2) BN, ReLU BiL, Conv(64, 3, 5, 1, 2) Tanh Output size : (3, 64, 64)	Input size : (8, 1, 1) Flatten FC(8, 256) ReLU FC(256, 1) Sigmoid  Output size : (1, )

Table 4: Encoder/Generator/Discriminator architectures for WAE-GAN for CelebA